

Reply to Comments from Reviewer #2

The referee comments are shown in blue.

The responses to the comments are shown in black.

The text included in the revised manuscript are shown in red.

Reviewer #2 comments:

This contribution details a region-specific, time-specific, and categorical evaluation of the meteorological and chemical forecasts from the offline-coupled GFSv15-CMAQv5.0.2 for the year 2019. This manuscript fits the scope of the journal of Geoscientific Model Development. However, the manuscript has important limitations. Although the paper is well organised and detailed, the results shown in the paper are, in my opinion, not sufficiently original and new to merit publication in this specific journal.

Response:

Thank you for your constructive comments. Please see below point-by-point responses.

Major comments

1. My main objection focuses on the fact that it is not clearly stated which elements of the study are genuinely new and original and which echo the findings of previous studies. I cannot clearly state either new methodology developed or any new results, apart from applying state-of-the-art validation methods to a new set of simulations with the GFSv15-CMAQv5.0.2.

Response:

Compared to previous publications in the literature, there are several new and unique aspects of our work.

Firstly, the historical and current NAQFC are based on the NAM-CMAQ system. However, the NAM has been no longer updated since March 2017, as NOAA Environmental Modeling Center (EMC) has transitioned to devote its full resources towards the development of an ensemble model based on the FV3 dynamical core. The FV3-based system will eventually replace all current NOAA National Centers for

Environmental Prediction (NCEP) mesoscale models used for forecasting. The next generation NAQFC will be based on the FV3. The NOAA National Weather Service (NWS) is currently coordinating an effort to inline a regional scale meteorological model basing on the same dynamic core as that in FV3GFS to be coupled with an atmospheric chemistry model partially based on CMAQ. The implementation of this inline system NAQFC is expected a few years in the future. Therefore, the system in this study is an important interim operational air quality forecasting capability the nation will have before the future inline system matures to be operational. The forecast skill of the interim FV3GFS-CMAQ system we analyzed in this work could also provide as a benchmark for the future inline system.

Secondly, while the FV3-based GFS v15 showed similar and even better performance compared to the previous GFS v14 (<https://www.emc.ncep.noaa.gov/users/meg/fv3gfs/>) and the NAM (https://hmt.noaa.gov/experiments/pdf/WPC-HMT_WWE_2019_Final_Report.pdf), air quality forecasting driven by FV3-based meteorology has not been fully evaluated. This is the first evaluation paper for the FV3GFS-CMAQ forecasting system, and demonstrates the ability of FV3-based GFSv15 to drive the air quality forecast.

Thirdly, an updated study is needed because previous journal publications for NAQFC were based on older versions of CMAQ (v4.6 or v4.7) and historical 2005 and 2011 NEIs (Kang et al., 2010a, 2010b; Garner et al., 2015; Bray et al., 2017; Huang et al., 2017; Lee et al., 2017; Pan et al., 2020). In our work, we implement 2014 NEI emission, GBBEPx fire emission, initial and lateral boundary conditions, and a CMAQ version closer to those in the current operational NAQFC thus allowing for the results and findings to be more directly transferable to the current NAQFC and inform further development of the future FV3-based NAQFC.

To further address the reviewer comments, we reorganize our manuscript in the revision. A section is added for detailed discussion of the processes and causes attributing to the biases. The analysis summarized in the evaluation section of the previous manuscript are incorporated and expanded for detailed analysis in the discussion section. Here we describe the major items added in the discussion section. We focus more on the impact of meteorology-chemistry relationship and the chemical processes on the model biases. Firstly, the meteorology-chemistry relationship is further quantified by adding the region-specific analysis of meteorological biases. Secondly, the discussion for major biases in O₃ and PM_{2.5} prediction is enhanced.

For O₃: The underpredicted O₃ during non-O₃ season was neglected in our previous manuscript. We include the temperature-O₃ relationship to quantify the impact of cold biases on the O₃ underprediction in the regions of 1, 3, 5, and 10, and on the overprediction in regions of 4 and 6 near the Gulf coast. While the cold biases could explain certain portion of the MDA8 O₃ underprediction, the dry deposition algorithm of O₃ to snow cover in CMAQ v5.0.2 is discussed for further contribution to the underprediction in the northern regions during winter and early spring. In addition to the temperature-O₃ relationships, the significant O₃ overprediction near the Gulf coast would be associated with the missing halogen chemistry and the overestimated emissions from oil-gas sources.

For PM_{2.5}: By this point, more and more observation datasets have been available. For example, most of the daily PM_{2.5} composition observations from Air Quality System (AQS) dataset are available for the entire 2019 now. We include them into further analysis and discussion for the manuscript. The contributions of the OC, sulfate, and dust compositions to the PM_{2.5} biases are quantified. Our findings are better supported by the additional evidence. These additional datasets and analysis help us to dig deeper into understanding the PM_{2.5} biases and their causes instead of referring previous studies to support our analysis in the original manuscript.

Our study supports the future development of the FV3GFS-CMAQ system and even the operational NAQFC, which has similar model settings and inputs with the system evaluated in this study. For example, by introducing the evaluation of PM_{2.5} compositions, we found the overestimation in dust compositions is one of the major sources for the PM_{2.5} overprediction in cooler months. We further conduct a sensitivity simulation by implementing an adjustment for suppressing the fugitive dust by snow cover, which was implemented in the operational NAQFC of the NAM-CMAQv5.0.2 system in 2020. We found that the suppression and adjustment could improve the model performance during cooler months. It further indicates the need of further development and improvement for the NAQFC.

2. Furthermore, objective (3) states that one of the aims of the manuscript is to "investigate underlying causes for the biases to provide a scientific basis for improving the model representations of chemical processes and developing science-based bias correction methods for O₃ and PM_{2.5} forecasts.". However, after reading the manuscript, I don't find any specific section of the manuscript devoted to understanding the

physicschemical processes causing over- or underproduction of air quality parameters (apart from some general comments). The authors should deepen into the processes leading the levels of air pollution so that objective (3) can really be achieved.

Response:

A discussion section is added for deepening the analysis of the underlying causes in three aspects: (1) Meteorology-chemistry relationships, (2) Causes for major biases in O₃ predictions, and (3) Causes for major biases in PM_{2.5} predictions. The revised text is listed below:

We further quantify the meteorology-chemistry relationships by conducting the region-specific evaluation of the meteorological variables. The regional performance for major variables is shown in Figure S9. The regional biases in T2 predictions show high correlation with the regional biases in MDA8 O₃. It indicates that the cold biases in the Midwest (including region 5) and the warm biases near the Gulf coast (including regions of 4 and 6) are important factors for the O₃ underprediction and overprediction in those regions, respectively. The O₃-temperature relationship was found (S. Sillman and Samson, 1995; Sillman, 1999). O₃ is expected to increase with increasing temperature within specific range of temperature (Bloomer et al., 2009; Shen et al., 2016). The surface MDA8 O₃-temperature relationship was found at approximately 3-6 ppb K⁻¹ in the eastern US (Rasmussen et al., 2012). According to such relationships, the biases in T2 predictions could explain large portion of the O₃ biases. Heavy convective precipitation and tropical cyclones occur more often in the southeastern US, which covers mainly regions 4 and 6. Therefore, the performance in precipitation predictions is worse in those two regions comparing to other regions as we have shown the model has relatively poor performance in capturing short-term heavy rains during summer seasons in section 3.1. Meanwhile, the performance in wind predictions in regions 4 and 6 is relatively poor. Such performance in the meteorological predictions is consistent with the mixed performance in PM_{2.5} prediction in regions 4 and 6. The discrepancy between simulated and observed meteorological variables, mainly in precipitations and wind, can be attributed to the poor temporal agreement shown as correlations of predicted PM_{2.5} in those two regions.

In the original submission, we focus on the significant O₃ overprediction near the Gulf coast during O₃-season. In the revised manuscript, we analyze likely causes

for additional model biases, e.g., underpredictions in the Northeast, Mid-Atlantic, Midwest, Mountainous states, and the Northwest (mainly corresponding to the regions 1, 3, 5, 8, and 9) during non-O₃ season as follows:

The O₃ concentration is underforecasted for the Northeast, Mid-Atlantic, Midwest, Mountainous states, and the Northwest (mainly corresponding to the regions 1, 3, 5, 8, and 9) during non-O₃ season. Large difference in dry deposition algorithms between CMAQ v5.0.2 and other common parameterizations was reported (Park et al., 2014; Wu et al., 2018). Large discrepancy between modeled dry deposition velocity of O₃ by CMAQ v5.0.2 and the observation during winter was shown and attributed to the deposition to snow surface. Improvement was indicated in revising the treatment of deposition to snow, vegetation, and bare ground in CMAQ v5.0.2. Lower deposition to snow was found to improve the consistency between the O₃ deposition modeled by CMAQ v5.0.2 and the observations. Therefore, the dry deposition module in v5.0.2 needs to be updated and improved for more accurate representation of low-moderate O₃ mixing ratios (Appel et al., 2020). For the cases in this study, the snow cover for the months of Jan and Apr in winter and spring is shown in Figure 7a and 7b. The underpredicted O₃ during non-O₃ season may be caused by the overestimated O₃ deposition to snow in the northern regions, corresponding to the previous regions 1, 3, 5, 8, and 9. The mixed effects of the temperature-O₃ relationship discussed above and the large deposition to snow contribute to the moderate O₃ underpredictions.

In the original manuscript, we conducted part of the analyses for the PM_{2.5} biases and the causes based on referring to some previous studies and several speculations. By introducing the additional AQS dataset into our evaluation, we can gain further insights into the specific reasons for the PM_{2.5} biases:

The variation in predicted PM_{2.5} composition between cooler and warmer months indicates that major seasonal biases are caused by multiple factors. We introduce the AQS dataset for evaluation of daily PM_{2.5} composition to provide additional insights into the specific reasons. Figure 9 shows the biases of the key PM_{2.5} composition for the cooler month of Jan and warmer month of Jul. While the overall mean biases of PM_{2.5} composition, including elemental carbon (EC), ammonium (NH₄⁺), and nitrate (NO₃⁻) are within $\pm 0.5 \mu\text{g m}^{-3}$ for all months of the year, the major biases in PM_{2.5} predictions are mostly contributed by organic carbon (OC), soil components (SOIL), and sulfate (SO₄²⁻). The soil components are estimated using the

Interagency Monitoring of Protected Visual Environments (IMPROVE) equation and specific constituents (Appel et al., 2013). During a cooler month, the significant overprediction in $PM_{2.5}$ is mainly attributed to the overprediction in OC and SOIL. During warmer months, the overprediction of SOIL and sulfate compensate for the overall underprediction in OC in v5.0.2, leading to the moderate $PM_{2.5}$ underprediction in the Southeast but slight overprediction in the Midwest, Mid-Atlantic, and the Northeast. These high $PM_{2.5}$ SOIL concentrations are consistent in spatial characteristics with large emissions of anthropogenic primary $PM_{2.5}$, and primary coarse PM in the Midwest, Northeast, and the Northwest (Fig. S6). The underprediction in $PM_{2.5}$ OC during summer compensate the overestimation in dust during cooler months, resulting in the overall biases with an annual NMB of 30%.

The large emissions of anthropogenic primary coarse PM, as well as the wind-blown dust are the major sources for predicted $PM_{2.5}$ SOIL components. Appel et al. (2013) indicated CMAQ overpredicted soil components in the eastern United States partially due to the anthropogenic fugitive dust and wind-blown dust emissions. The overprediction in $PM_{2.5}$ soil compositions by our forecast system could be mainly attributed to the overestimation of the anthropogenic fugitive dust emission because the meteorological conditions were not included in processing the anthropogenic fugitive dust sector. The dust-related components of aluminum, calcium, iron, titanium, silicon, and coarse mode particles are overestimated in the regions with snow and precipitation, especially during winter, early spring, and late autumn with snow cover in the north, which contributes to the $PM_{2.5}$ overprediction, with more significant temporal-spatial pattern in the northern U.S. during cooler months.

In CMAQ v5.0.2, the primary organic aerosol (POA) is processed as non-volatile. The emissions of semivolatile and intermediate volatility organic compounds (S/IVOCs) and their contributions to the secondary organic aerosol (SOA) are not accounted for in the aerosol module. In the recent versions of CMAQ, two approaches linked to POA sources have been implemented. One introduces semi-volatile partitioning and gas-phase oxidation of POA emissions. The other (called pcSOA) accounts for multiple missing sources of anthropogenic SOA formation, including potential missing oxidation pathways and emissions of IVOCs. These two improvements lead to increased organic carbon concentrations in summer but decreased level in winter. The changes vary by season as a result of differences in volatility (as

dictated by temperature and boundary layer height) and reaction rate between winter and summer. Therefore, the missing S/IVOCs and related SOA chemistry in v5.0.2 are key reasons for the OC overprediction and underprediction during cooler and warmer months, respectively.

The detailed analysis for the underlying cause by the uncertainty and biases in emissions of anthropogenic fugitive dust will be discussed further in the discussion section and in the response to the “other comment #2 by Reviewer 2”.

In general, we reorganize the evaluation and the discussion sections. The discussion of relationships between the meteorological factors and the chemical performance are further enhanced. The physical drivers are linked with the biases in air pollutant predictions more deeply. The discussion of the model chemistry is expanded in more detail.

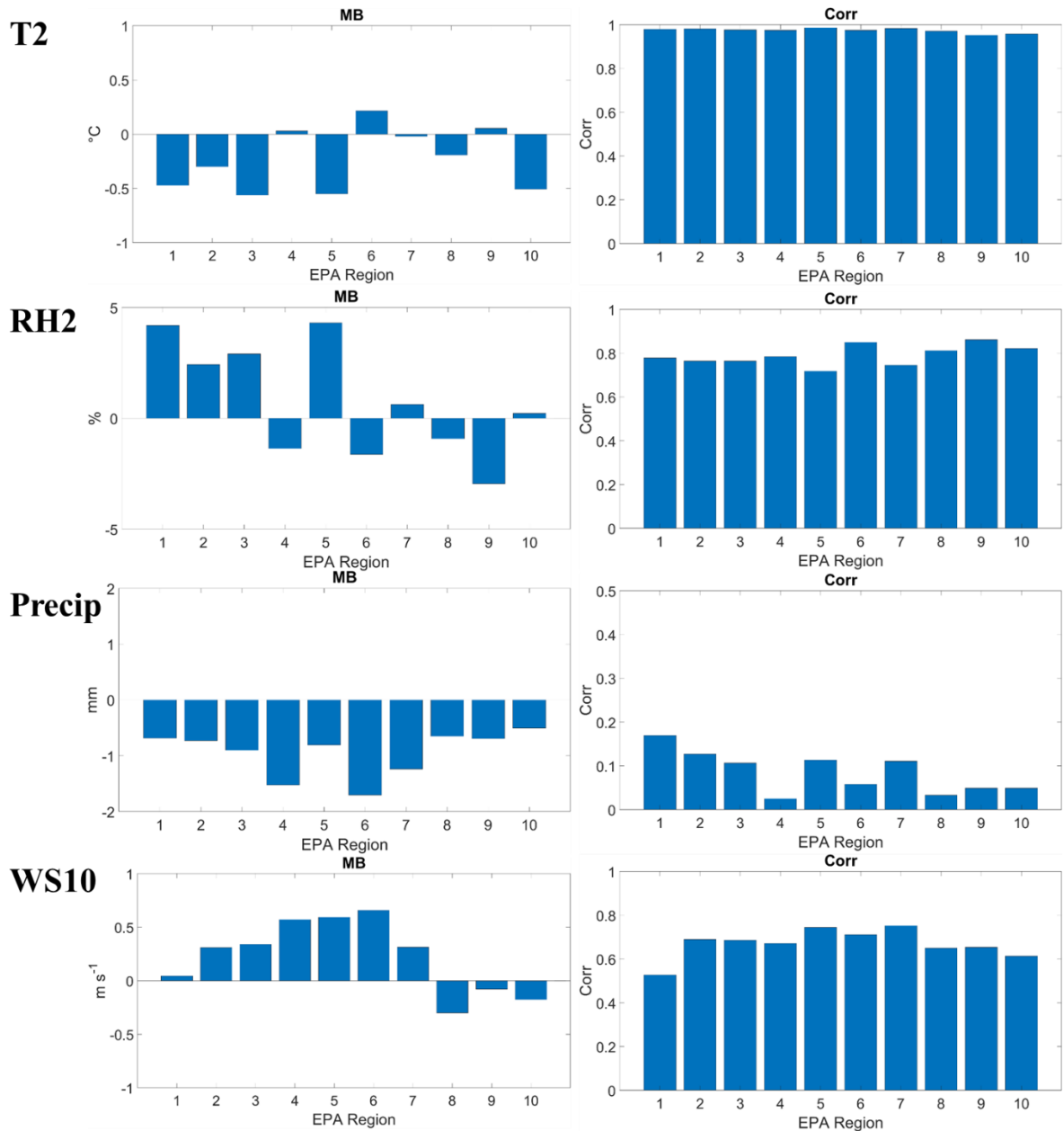


Figure S9. Annual performance for 10 regions in predicting meteorological variables of temperature at 2-m (T2), relative humidity at 2-m (RH2), precipitation, and wind speed at 10-m (WS10)

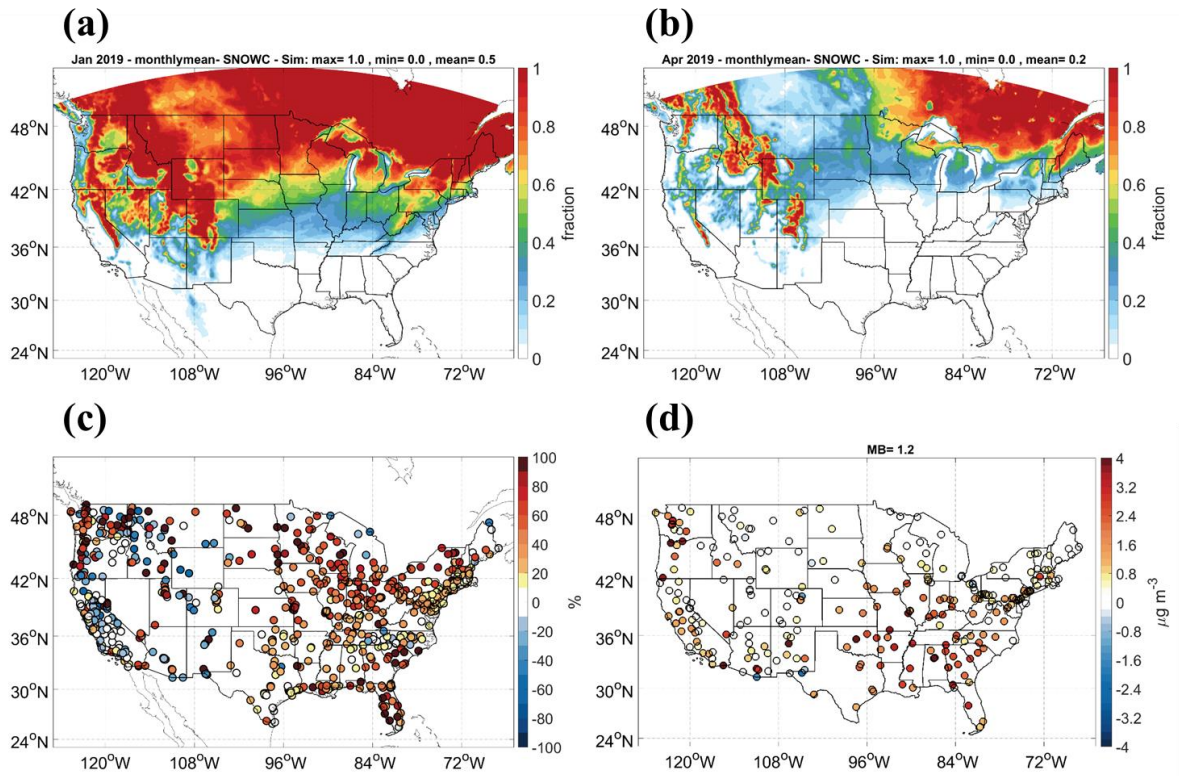


Figure 7. The predicted average snow cover for (a) Jan and (b) Apr. (c) The difference in NMBs of $PM_{2.5}$ by adjusting anthropogenic fugitive dust emission for Jan. Positive values stand for improvement in biases with NMBs closer to 0. (d) MBs in $PM_{2.5}$ soil composition with adjustment of AFD emission for Jan.

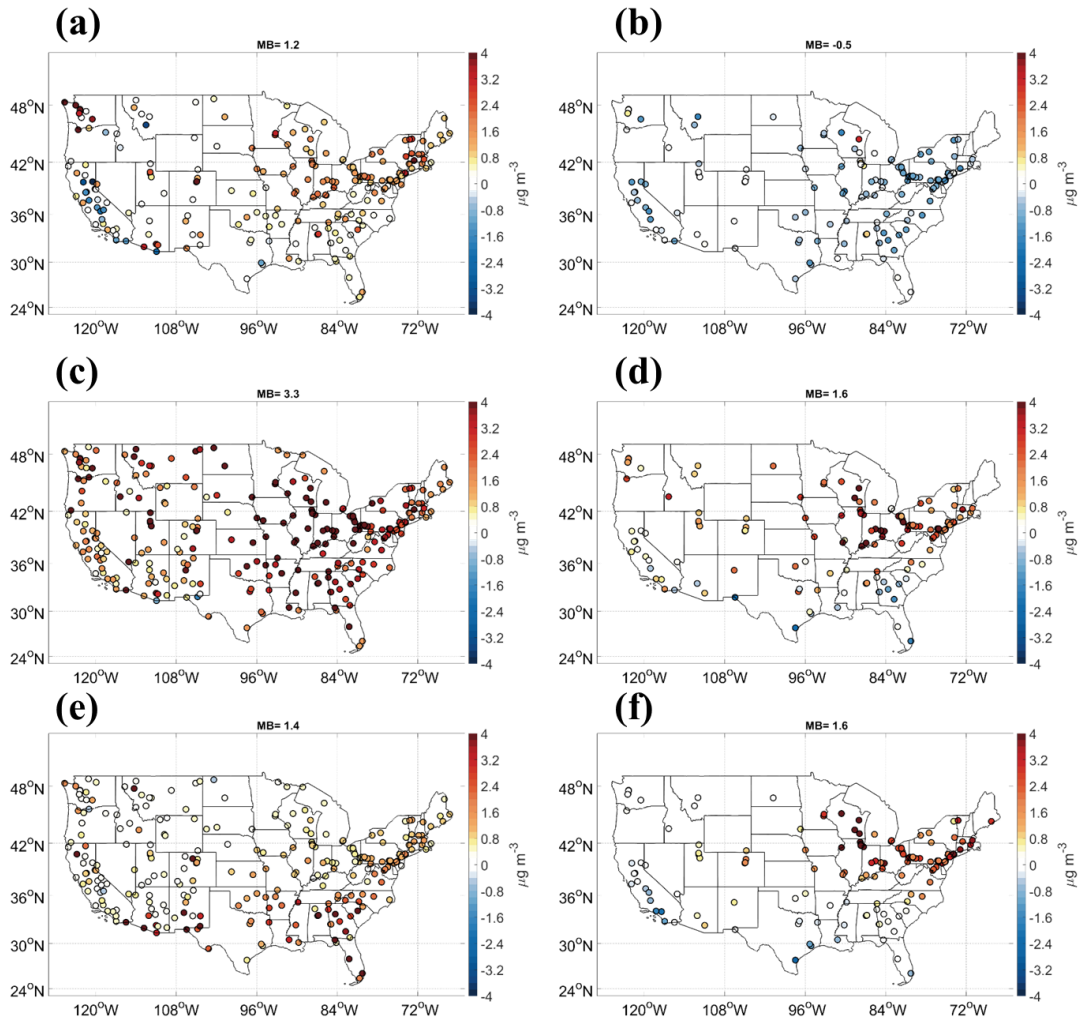


Figure 9. Mean biases in PM_{2.5}: (a) OC for Jan, (b) OC for Jul, (c) SOIL for Jan, (d) SOIL for Jul, (e) sulfate for Jan, and (f) sulfate for Jul

3. Despite the large number of statistical figures presented, I have the very personal opinion that the authors do not take the advantage of the compiled information to point to the specific causes for model biases.

Response:

We agree that some of the presented figures and materials were not well interpreted to informative messages to the readers. For example, the evaluation of monthly accumulated precipitation is shown in three figures (Figures S2 to S4 in the original submission). But the main purpose of those three figures in the original manuscript was only to indicate that the GFSv15-CMAQv5.0.2 system has better agreement in the spatial characteristics than the temporal variations. We revise the presentation of the evaluation of accumulated precipitation as Figure S3. The

performance and the spatial patterns of the monthly accumulated precipitation during four seasons are more straightforward. Meanwhile, the findings from the evaluation section 3 are further discussed in section 4 in details. We closely link the meteorological drivers to the analysis for chemical biases as we indicated in the response of comment #2. The cold biases shown in Figure S2 are further discussed in section 4.1. Additional information and analysis are added in section 4.2 to address the major biases in O₃ predictions shown in Figures 2 and 6. Furthermore, the information and description for Figure 7 in the original manuscript are revised and enhanced in section 4.3 by introducing the evaluation of PM_{2.5} against AQS dataset. The original Figure 9 is further supported by the additional analysis of seasonal variation of PM_{2.5} composition and the diurnal emissions as indicated in the response to major comment #2 and other comment #2.

In general, we reorganize some figures to make them more concise. We deepen the analysis by providing additional supporting material for several speculations in the revised discussion section.

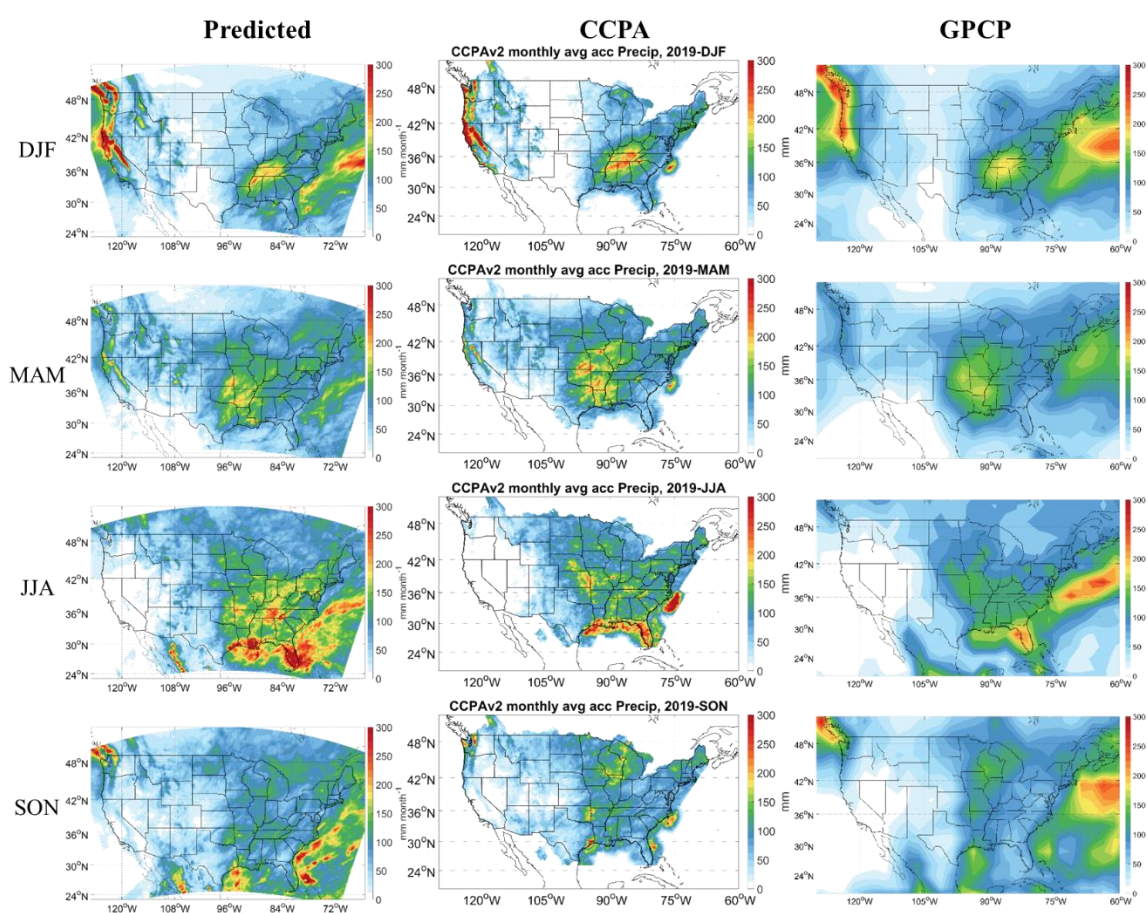


Figure S3. Monthly accumulated precipitation for four seasons by the GFSv15-CMAQv5.0.2 prediction, CCPA observation, and GPCP observation

Other comments

1. The authors should compare the skills of the model (categorical evaluation) with other published model studies, in order to have a flavor of the behavior of the model when compared to other forecasting systems worldwide.

Response:

Major RT-AQF systems over the world were comprehensively reviewed in (Zhang et al., 2012a, 2012b). Here we include a comparison with the more recent air quality forecasting studies from Canada (Moran et al., 2018; Russell et al., 2019), Europe (Struzewska et al., 2016; D'Allura et al., 2018; Podrascanin, 2019; Stortini et al., 2020), East Asia (Lyu et al., 2017; Zhou et al., 2017; Peng et al., 2018; Ha et al., 2020), and CONUS (Kang et al., 2010; Zhang et al., 2016; Lee et al., 2017). We summarize the performance in these studies in a table in the supplementary material. As for the categorical performance, the air quality standards vary in different regions (Oliveri Conti et al., 2017). For example, National Ambient Air Quality Standards (NAAQSs), the Ambient Air Quality and Cleaner Air for Europe (CAFE) Directive (2008/50/EC), and the national ambient air quality standard (GB 3095-2012) are set up by U.S., Europe, and China, respectively. Therefore, the definition of the categorical metrics may vary between regions even with the same metric name. Their categorical performance are discussed specifically in the revised text:

Table S3 summarizes air quality forecasting skills reported in the literature along with that from this work. For those studies with data assimilation in air quality forecasting, the performance from the raw results without data assimilation are presented. The performance in predicting O₃ and PM varies largely between model systems. The discrete and categorical performance in O₃ prediction is not significantly better than that in PM prediction. O₃ tends to be slightly overpredicted in an annual base or for the warmer months. The annual NMB and Corr for O₃ over the North America are 1.4% and 0.76 for 2010 in Moran et al. (2018), while they are 1.0% and 0.73 in this study. However, the performance in PM_{2.5} prediction varies largely from our study. The PM_{2.5} for warmer months were moderately overpredicted in Russel et al. (2019), with the MBs ranging from 3.2 to 5.5 μg m⁻³. The categorical performance of GFSv15-CMAQv5.0.2 in predicting MDA8 O₃ is similar with that of the previous

NAQFC (Kang et al., 2010), in which the FAR and H are ~68 % and ~31% for “Unhealthy for Sensitive Groups”, and the H is ~47% for “Moderate” category, respectively. The H for PM_{2.5} also decreased largely from ~46% for “Moderate” to ~21% for “Unhealthy for Sensitive Groups” category, and the FAR was over 90% for the “Unhealthy for Sensitive Groups” category in Kang et al. (2010). The overpredicted PM_{2.5} was also found when using the historical 2005 NEI in forecast for Jan 2015 (Lee et al., 2017). The performance was improved by updates of 2011 NEI and real-time dust and wildfire emissions. It indicates the needs of improving our emission inventory. As for the categorical performance in regions other than CONUS, the air quality standards vary (Oliveri Conti et al., 2017). For example, National Ambient Air Quality Standards (NAAQSs), the Ambient Air Quality and Cleaner Air for Europe (CAFE) Directive (2008/50/EC), and the national ambient air quality standard (GB 3095-2012) are set up by U.S., Europe, and China, respectively. Metrics also vary between studies. The primary forecasting products are O₃ and PM₁₀ from some forecasting systems instead of O₃ and PM_{2.5} in this study. The threshold for categorical evaluation of O₃ used in D’Allura et al (2018) was 83.0 µg m⁻³. The applied metrics of the False Alarm Ratio and Probability of Detection (POD) were defined the same as the FAR and H used in our study. The FAR and POD were 36.14% and 71.16%, respectively. The categorical evaluation of PM_{2.5} in Ha et al. (2020) was applied for four categories: (1) 0-15 µg m⁻³, (2) 16-50 µg m⁻³, (3) 51-100 µg m⁻³, and (4) >100 µg m⁻³. The overall FAR and Detection Rate for four categories are 59.0% and 36.1%, respectively. Although the metrics of FAR and Detection Rate were defined for four categories, rather than every single category as for this study, the categorical performance is comparable with our results. In general, the discrete and categorical performance of O₃ forecast in this study is comparable that of the air quality forecasting systems in many regions of the world. However, the PM forecasts vary largely between studies. While our GFSv15-CMAQv5.0.2 system shows consistent performance with the systems covering CONUS, the high FAR and low H for “Unhealthy for Sensitive Groups” category with higher thresholds indicate that the categorical performance could be further improved by addressing the significant overprediction during cooler months in this study.

Table S3. Summary of forecasting skills of air quality forecasting systems

Reference	Region	Model System	Period	Pollutant	Performance
Moran et al., 2018	Canada/North America	GEM-MACH	2010	O ₃	NMB=1.4%, R=0.76
				PM _{2.5}	NMB=-0.6%, R=0.58
Russell et al., 2019	Canada	GEM-MACH	Aug-Sept 2013	O ₃	MB=5.7 to 10.9 ppb, RMSE=9.7 to 16.0 ppb, Corr=0.50 to 0.74
				PM _{2.5}	MB=3.2 to 5.5 $\mu\text{g m}^{-3}$, RMSE=5.7 to 8.8 $\mu\text{g m}^{-3}$, Corr=0.20 to 0.47
Struzewska et al., 2016	Poland	GEM-AQ	Nov 2011 to Sep 2013	O ₃	MB=12.8 to 25.6 $\mu\text{g m}^{-3}$, RMSE=24.6 to 28.7 $\mu\text{g m}^{-3}$, Corr=0.48 to 0.62
				PM _{2.5}	MB=-9.6 to -1.86 $\mu\text{g m}^{-3}$, RMSE=24.8 to 34.1 $\mu\text{g m}^{-3}$, Corr=0.48 to 0.58
D'Allura et al., 2018	Italy	WRF/RAMS-FARM	2015	O ₃	FAR=36.1%, POD= 71.2%, threshold=83 $\mu\text{g m}^{-3}$
				PM ₁₀	FAR=20.0%, POD= 27.3%, threshold=33 $\mu\text{g m}^{-3}$
Podrascanin, 2019	Serbia	WRF/Chem	August 2016	O ₃	MB=1.6 to 9.3 $\mu\text{g m}^{-3}$, NMB=3.0 to 17.2%, Corr=0.45 to 0.50
				PM ₁₀	MB=-15.2 to -14.3 $\mu\text{g m}^{-3}$, NMB=-74.0 to -56.1%, Corr=-0.01 to 0.18
Stortini et al., 2020	Italy	CHIMERE	October 2019	O ₃	MB=11.0 to 16.9 $\mu\text{g m}^{-3}$, RMSE=19.3 to 28.0 $\mu\text{g m}^{-3}$, Corr=0.63 to 0.78
				PM ₁₀	MB=-8.2 to -4.9 $\mu\text{g m}^{-3}$, RMSE=11.4 to 13.0 $\mu\text{g m}^{-3}$, Corr=0.72 to 0.76, FAR=43-44%, POD= 6-22%
Lyu et al., 2017	China	WRF-CMAQ	2014-2016	PM _{2.5}	NME=49%, RMSE=32.2 $\mu\text{g m}^{-3}$, R ² =0.46
Zhou et al., 2017	China	WRF/Chem	2014-2015	MDA8	MB=18.9 ppb, NMB=77%, RMSE=27.9 ppb, Corr=0.63
				PM _{2.5}	MB=-12.0 $\mu\text{g m}^{-3}$, NMB=-9%, RMSE=35.8 $\mu\text{g m}^{-3}$, Corr=0.67
Peng et al., 2018	China	WRF/Chem	6 to 16 October 2014	O ₃	MB=-31.0 $\mu\text{g m}^{-3}$, RMSE=50.8 $\mu\text{g m}^{-3}$, Corr=0.46
				PM _{2.5}	MB=-34.1 $\mu\text{g m}^{-3}$, RMSE=92.1 $\mu\text{g m}^{-3}$, Corr=0.74
Ha et al., 2020	Korea	WRF/Chem	May 2016	PM _{2.5}	MAE=12.8, FAR=59.0%, Overall_Accuracy=59.7%, High_Pollution_Accuracy=35.6%
Kang et al., 2010	CONUS	NAM-CMAQ	2008	MDA8	HR=~47% for cat 2, ~31% for cat 3, FAR=~68% at cat 3
				PM _{2.5}	HR=~46% for cat 2, ~21% for cat 3, FAR=~91% at cat 3
Zhang et al., 2016	Southeastern US	WRF/Chem-MADRID	2012-2014 (May-September)	MDA8	NMB= 0.0 to 17.0 %, NME= 22.0 to 27.0 %, Corr= 0.5
				PM _{2.5}	NMB= -4.0 to 15.0 %, NME= 36.0 to 40.0 %, Corr= 0.3 to 0.4

			2012-2014 (December- February)	MDA8 O ₃	NMB= -17.7 to -9.1 %, NME= 19.6 to 24.6 %, Corr= 0.0 to 0.2
				PM _{2.5}	NMB= 0.8 to 8.3 %, NME= 42.6 to 47.2 %, Corr= 0.3 to 0.4
Lee et al., 2017	CONUS	NAM- CMAQ	May and Jul 2014	PM _{2.5}	MB=-2.7 to -1.6 µg m ⁻³ , NME=-35 to -20%, RMSE=4.5 to 8.8 µg m ⁻³ , Corr=0.22 to 0.33
			Jan 2015	PM _{2.5}	MB=1.3 to 3.7 µg m ⁻³ , NME=13 to 38%, RMSE=6.5 to 12.6 µg m ⁻³ , Corr=-0.37 to 0.38
This work	CONUS	FV3GFS- CMAQ	2019	MDA8 O ₃	MB=0.4 ppb, NMB=1.0%, Corr=0.73. FAR=41.4 % and HR=45.8% at cat 2.
				PM _{2.5}	MB=2.3 µg m ⁻³ , NMB=30.0%, Corr=0.41. FAR=70.3% and HR=57.6% for cat 2.

FAR: False Alarm Ratio; POD: Probability of Detection; HR: Hit Rate.

2. Emissions are really important in forecasting system; however, this manuscript lacks information about the emissions used (time series, spatial patterns, seasonal behavior, etc). The authors should explain in a higher degree of detail how emissions are considered and implemented in the forecasting system.

Response:

We expand the introduction of how we prepare and implement the emission into the GFSv15-CMAQv5.0.2 system in the methodology section 2:

The anthropogenic emissions from area, mobile, and point sources in National Emissions Inventory of year 2014 version 2 (NEI 2014v2) are processed by the Sparse Matrix Operator Kernel Emissions (SMOKE) modeling system. The onroad mobile sources include all emissions from motor vehicles that operate on roadways such as passenger cars, motorcycles, minivans, sport-utility vehicles, light-duty trucks, heavy-duty trucks, and buses. Onroad mobile source emissions are processed using emission factors output from the Motor Vehicle Emissions Simulator (MOVES). SMOKE uses a combination of vehicle activity data, emission factors from MOVES, meteorology data, and temporal allocation information to estimate hourly, gridded onroad emissions. The nonroad, agriculture, anthropogenic fugitive dust, non-elevated oil-gas, residential wood combustion, and other sectors are included in the area sources. The sectors of airports, commercial marine vessel (CMV), electric generating units (pt_egu), point

sources related to oil and gas production (pt_oilgas), point sources that are not EGUs nor related to oil and gas (ptnonipm), and point sources outside US (pt_other) are included in the point sources. The sulfur dioxide (SO₂) and nitrogen oxide (NO_x) from point sources in NEI 2005 are projected to year 2019 following the methods used in Tang et al. (2015, 2017). The biomass burning emission inventory from the Blended Global Biomass Burning Emissions Product system (GBBEPx) (Zhang et al., 2019b) is implemented for the forecast of forest fires. The GBBEPx fire emission is treated as one type of point source. Its heat flux is derived from satellite retrieved fire radiative power (FRP) to drive fire plume rise. The GBBEPx is a near real time fire dataset. The fire emission implemented in the current forecast cycle comes from the historical fire observation, typically 1-2 day behind. In this system, we use landuse information to classify fires into forest fire and other burning such as agriculture burning. We assume only forest fire can last longer than 24 hours. We assume the forest fire emission will continue on day 2 and beyond. Other types of fires will be dropped as we assume few of them could continue beyond day 2. The plume rise of the point source is driven by the meteorology and allocated to the 35 elevated layers in GFSv15-CMAQv5.0.2 system by the PREMAQ preprocessing system. Biogenic emissions are calculated inline by Biogenic Emission Inventory System (BEIS) version 3.14 (Schwede et al., 2005). Sea-salt emission is parameterized within CMAQ v5.0.2. While the deposition velocities are calculated inline, the fertilizer ammonia bi-directional flux for in-line emissions and deposition velocities is turned off.

The impact of the emissions on the biases in O₃ prediction is added in the revised text:

In addition to the impact of meteorological biases and missing halogen chemistry on the O₃ overprediction near Gulf coast, the overestimated VOC emission could increase O₃ biases. The anthropogenic VOCs emissions continuously decrease from historical NEIs to 2016 NEI (<http://views.cira.colostate.edu/wiki/wiki/10202/inventory-collaborative-2016v1-emissions-modeling-platform>). We compare the VOCs emissions between 2016 NEI and the emissions used in this study. Figure S10 shows the difference in the elevated source of pt_oilgas. The Gulf coast is impacted by the oil and gas sector due to the oil and gas fields, and the exploration activity near it. By comparing the 2016 NEI to the current emissions we used in the system, we found that the overestimation of the VOCs emissions could be one aspect to the O₃ overprediction near the Gulf Coast. We only

project the SO₂ and NO_x from 2005 NEI to 2019 and we do not project the VOCs for the elevated sources. The monthly VOCs emissions from pt_oilgas sector for July in regions 4 and 6 are 2876.0 tons month⁻¹, while they are 2497.0 tons month⁻¹ in 2016 NEI. The reduction mainly locates along the coastline, where the significant overprediction takes place. It indicates the complicated effect of meteorological biases, missing gas-phase chemistry, and the overestimation of emissions on the O₃ prediction in such area.

While the diurnal characteristics of the emissions with the revised Figure S11 are added in the revised manuscript to understand the diurnal PM_{2.5} biases, additional analyses are conducted for specific issues below:

During cooler months, the significantly overpredicted PM_{2.5} is mainly attributed to the emission of anthropogenic fugitive dust. In reality, the meteorological conditions could largely impact the amount and characteristics of anthropogenic fugitive dust. For example, the snow cover and the soil moisture are important factors in calculating the dust emissions in SMOKE. However, the anthropogenic fugitive dust implemented in this GFSv15-CMAQv5.0.2 system was not adjusted by the precipitation and snow cover. The large emissions of anthropogenic primary coarse PM, as well as the wind-blown dust are the major sources for predicted PM_{2.5} SOIL components. Appel et al. (2013) indicated CMAQ overpredicted soil components in the eastern United States partially due to the anthropogenic fugitive dust and wind-blown dust emissions. The overprediction in PM_{2.5} soil compositions by our forecast system could be mainly attributed to the overestimation of the anthropogenic fugitive dust emission because the meteorological conditions were not included in processing the anthropogenic fugitive dust sector. The dust-related components of aluminum, calcium, iron, titanium, silicon, and coarse mode particles are overestimated in the regions with snow and precipitation, especially during winter, early spring, and late autumn with snow cover in the north. Thus, it contributes to the PM_{2.5} overprediction, with more significantly temporal-spatial pattern in the north U.S. during cooler months.

An adjustment of precipitation and snow cover for fugitive dust was implemented in the operational NAQFC. The dust-related PM emissions will be clean up using a factor of 0.01 when the snow cover is higher than 25% or the hourly precipitation is higher than 0.1 mm hr⁻¹ before they are used as input for CMAQ v5.0.2 forecast. We conduct a sensitivity simulation for Jan 2019 using the GFSv15-CMAQv5.0.2 system with the adjustment implemented in the operational NAQFC.

Figure 7c shows the PM_{2.5} overprediction in the northern regions 1, 2, 5, and 10 during Jan is largely improved corresponding to the spatial-temporal characteristics of snow cover. The monthly MB and NMB for Jan improves from 5.5 $\mu\text{g m}^{-3}$ and 66.9% to 2.1 $\mu\text{g m}^{-3}$ and 24.0%, respectively. The improvement is mainly attributed to the decrease in overpredictions in PM_{2.5} soil components, with MBs decreased from 3.3 $\mu\text{g m}^{-3}$ to 1.2 $\mu\text{g m}^{-3}$ for Jan (Fig. 7d). The overprediction in the Northeast and Northwest during spring is expected to be improved by the suppression of the fugitive dust by the snow during early spring. This indicates the importance of including the meteorological forecast in processing the emission of anthropogenic fugitive dust. It should be calculated inline or be adjusted by the meteorological forecast.

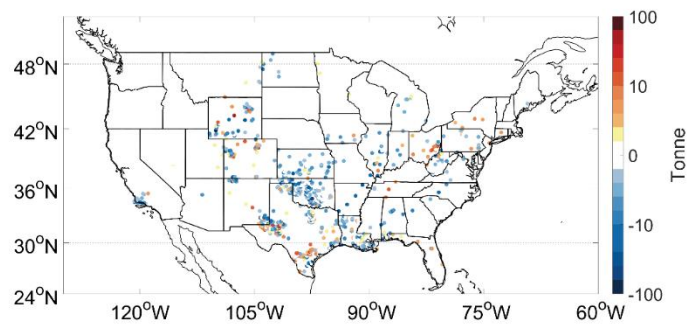


Figure S10. Difference of VOC emissions from pt_oilgas sector in 2016 NEI comparing to the emissions used in this study

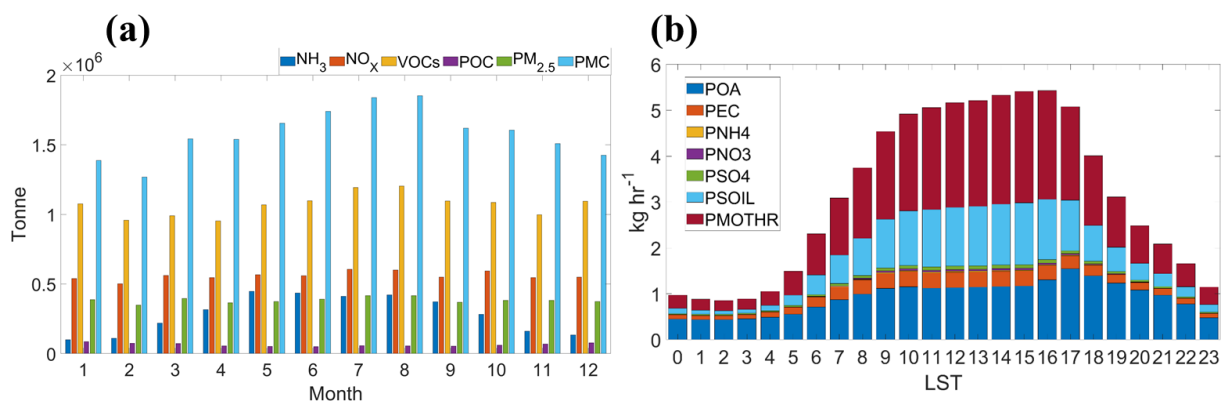


Figure S11. (a) Monthly variation of domain-wide surface emissions, and (b) diurnal emissions of fine mode PM

Reference

- Appel, K. W., Bash, J., Fahey, K., Foley, K., Gilliam, R., Hogrefe, C., et al. (2020). The Community Multiscale Air Quality (CMAQ) Model Versions 5.3 and 5.3.1: System Updates and Evaluation. *Geoscientific Model Development Discussions*, 1–41. <https://doi.org/10.5194/gmd-2020-345>
- Bray, C. D., Battye, W., Aneja, V. P., Tong, D., Lee, P., Tang, Y., & Nowak, J. B. (2017). Evaluating ammonia (NH₃) predictions in the NOAA National Air Quality Forecast Capability (NAQFC) using in-situ aircraft and satellite measurements from the CalNex2010 campaign. *Atmospheric Environment*, 163, 65–76. <https://doi.org/10.1016/j.atmosenv.2017.05.032>
- D'Allura, A., Costa, M. P., & Silibello, C. (2018). Qualearia: European and national scale air quality forecast system performance evaluation. *International Journal of Environment and Pollution*, 64(1–3), 110–124. <https://doi.org/10.1504/IJEP.2018.099152>
- Garner, G. G., Thompson, A. M., Lee, P., & Martins, D. K. (2015). Evaluation of NAQFC model performance in forecasting surface ozone during the 2011 DISCOVER-AQ campaign. *Journal of Atmospheric Chemistry*, 72(3–4), 483–501. <https://doi.org/10.1007/s10874-013-9251-z>
- Ha, S., Liu, Z., Sun, W., Lee, Y., & Chang, L. (2020). Improving air quality forecasting with the assimilation of GOCI aerosol optical depth (AOD) retrievals during the KORUS-AQ period. *Atmospheric Chemistry and Physics*, 20(10), 6015–6036. <https://doi.org/10.5194/acp-20-6015-2020>
- Huang, J., McQueen, J., Wilczak, J., Djalalova, I., Stajner, I., Shafran, P., et al. (2017). Improving NOAA NAQFC PM 2.5 Predictions with a Bias Correction Approach. *Weather and Forecasting*, 32(2), 407–421. <https://doi.org/10.1175/WAF-D-16-0118.1>
- Kang, D., Mathur, R., & Trivikrama Rao, S. (2010). Assessment of bias-adjusted PM_{2.5} air quality forecasts over the continental United States during 2007. *Geoscientific Model Development*, 3(1), 309–320. <https://doi.org/10.5194/gmd-3-309-2010>
- Kang, Daiwen, Mathur, R., & Trivikrama Rao, S. (2010). Real-time bias-adjusted O₃ and PM_{2.5} air quality index forecasts and their performance evaluations over the continental United States. *Atmospheric Environment*, 44(18), 2203–2212. <https://doi.org/10.1016/j.atmosenv.2010.03.017>
- Lee, P., McQueen, J., Stajner, I., Huang, J., Pan, L., Tong, D., et al. (2017). NAQFC Developmental Forecast Guidance for Fine Particulate Matter (PM 2.5). *Weather and Forecasting*, 32(1), 343–360. <https://doi.org/10.1175/waf-d-15-0163.1>
- Lyu, B., Zhang, Y., & Hu, Y. (2017). Improving PM_{2.5} Air Quality Model Forecasts in China Using a Bias-Correction Framework. *Atmosphere*, 8(12), 147. <https://doi.org/10.3390/atmos8080147>
- Moran, M. D., Lupu, A., Zhang, J., Savic-Jovcic, V., & Gravel, S. (2018). A comprehensive performance evaluation of the next generation of the canadian operational regional air quality deterministic prediction system. In *Springer Proceedings in Complexity* (pp. 75–81). Springer. <https://doi.org/10.1007/978-3->

- Oliveri Conti, G., Heibati, B., Kloog, I., Fiore, M., & Ferrante, M. (2017). A review of AirQ Models and their applications for forecasting the air pollution health outcomes. *Environmental Science and Pollution Research*, *24*(7), 6426–6445. <https://doi.org/10.1007/s11356-016-8180-1>
- Pan, L., Kim, H., Lee, P., Saylor, R., Tang, Y., Tong, D., et al. (2020). Evaluating a fire smoke simulation algorithm in the National Air Quality Forecast Capability (NAQFC) by using multiple observation data sets during the Southeast Nexus (SENEX) field campaign. *Geoscientific Model Development*, *13*(5), 2169–2184. <https://doi.org/10.5194/gmd-13-2169-2020>
- Park, R. J., Hong, S. K., Kwon, H.-A., Kim, S., Guenther, A., Woo, J.-H., & Loughner, C. P. (2014). An evaluation of ozone dry deposition simulations in East Asia. *Atmospheric Chemistry and Physics*, *14*(15), 7929–7940. <https://doi.org/10.5194/acp-14-7929-2014>
- Peng, Z., Lei, L., Liu, Z., Sun, J., Ding, A., Ban, J., et al. (2018). The impact of multi-species surface chemical observation assimilation on air quality forecasts in China. *Atmospheric Chemistry and Physics*, *18*(23), 17387–17404. <https://doi.org/10.5194/acp-18-17387-2018>
- Podrascanin, Z. (2019). Setting-up a Real-Time Air Quality Forecasting system for Serbia: a WRF-Chem feasibility study with different horizontal resolutions and emission inventories. *Environmental Science and Pollution Research*, *26*(17), 17066–17079. <https://doi.org/10.1007/s11356-019-05140-y>
- Russell, M., Hakami, A., Makar, P. A., Akingunola, A., Zhang, J., Moran, M. D., & Zheng, Q. (2019). An evaluation of the efficacy of very high resolution air-quality modelling over the Athabasca oil sands region, Alberta, Canada. *Atmospheric Chemistry and Physics*, *19*(7), 4393–4417. <https://doi.org/10.5194/acp-19-4393-2019>
- Spiridonov, V., Jakimovski, B., Spiridonova, I., & Pereira, G. (2019). Development of air quality forecasting system in Macedonia, based on WRF-Chem model. *Air Quality, Atmosphere and Health*, *12*(7), 825–836. <https://doi.org/10.1007/s11869-019-00698-5>
- Stortini, M., Arvani, B., & Deserti, M. (2020). Operational forecast and daily assessment of the air quality in Italy: A copernicus-CAMS downstream service. *Atmosphere*, *11*(5), 447. <https://doi.org/10.3390/ATMOS11050447>
- Struzewska, J., Kaminski, J. W., & Jefimow, M. (2016). Application of model output statistics to the GEM-AQ high resolution air quality forecast. *Atmospheric Research*, *181*, 186–199. <https://doi.org/10.1016/j.atmosres.2016.06.012>
- Tang, Y., Chai, T., Pan, L., Lee, P., Tong, D., Kim, H.-C., & Chen, W. (2015). Using optimal interpolation to assimilate surface measurements and satellite AOD for ozone and PM_{2.5}: A case study for July 2011. *Journal of the Air & Waste Management Association*, *65*(10), 1206–1216. <https://doi.org/10.1080/10962247.2015.1062439>
- Tang, Y., Pagowski, M., Chai, T., Pan, L., Lee, P., Baker, B., et al. (2017). A case study of aerosol data assimilation with the Community Multi-scale Air Quality Model

over the contiguous United States using 3D-Var and optimal interpolation methods. *Geoscientific Model Development*, 10(12), 4743–4758. <https://doi.org/10.5194/gmd-10-4743-2017>

Wu, Z., Schwede, D. B., Vet, R., Walker, J. T., Shaw, M., Staebler, R., & Zhang, L. (2018). Evaluation and Intercomparison of Five North American Dry Deposition Algorithms at a Mixed Forest Site. *Journal of Advances in Modeling Earth Systems*, 10(7), 1571–1586. <https://doi.org/10.1029/2017MS001231>

Zhang, Y., Bocquet, M., Mallet, V., Seigneur, C., & Baklanov, A. (2012a, December 1). Real-time air quality forecasting, part I: History, techniques, and current status. *Atmospheric Environment*. Pergamon. <https://doi.org/10.1016/j.atmosenv.2012.06.031>

Zhang, Y., Bocquet, M., Mallet, V., Seigneur, C., & Baklanov, A. (2012b, December 1). Real-time air quality forecasting, Part II: State of the science, current research needs, and future prospects. *Atmospheric Environment*. Pergamon. <https://doi.org/10.1016/j.atmosenv.2012.02.041>

Zhang, Y., Hong, C., Yahya, K., Li, Q., Zhang, Q., & He, K. (2016). Comprehensive evaluation of multi-year real-time air quality forecasting using an online-coupled meteorology-chemistry model over southeastern United States. *Atmospheric Environment*, 138, 162–182. <https://doi.org/10.1016/j.atmosenv.2016.05.006>

Zhou, G., Xu, J., Xie, Y., Chang, L., Gao, W., Gu, Y., & Zhou, J. (2017). Numerical air quality forecasting over eastern China: An operational application of WRF-Chem. *Atmospheric Environment*, 153, 94–108. <https://doi.org/10.1016/j.atmosenv.2017.01.020>